

Application of K-Nearest Neighbors to Fraudulent Credit Card Transaction Classification

1 Introduction

The K-Nearest Neighbors (K-NN) model stands out for its simplicity and effectiveness among the algorithms we used to classify the fraudulent credit card transactions dataset. This section of the report examines the application of the K-NN model to the dataset, covering data preparation, feature selection, model training, and evaluation. By the end of this section, readers will gain a comprehensive understanding of how the K-NN algorithm was implemented in our project, the challenges encountered, and the outcomes achieved.

Despite its simplicity, the K-NN algorithm can achieve high performance when properly implemented and applied to the right context. For classification tasks, the algorithm assigns class labels by determining which data points in the training set are most similar to each data point in the test set. The class label is chosen from the k most similar "neighbors" through a majority vote. In the context of K-NN, similarity is measured using a distance metric, with common metrics including Euclidean, Manhattan, and Minkowski distances.

In this section, we will explore the steps involved in implementing the K-NN model for our dataset. We will start with data preparation, ensuring the data is clean and normalized. Next, we will discuss feature selection to identify the most relevant features for classification. The model training phase will cover the selection of the optimal k value, followed by an evaluation of the model's performance using cross-validation. Finally, we will reflect on the results obtained, describing key insights and any challenges we faced during the implementation process.

This report examines each phase of the K-NN model application to provide a clear and detailed account of how we used the algorithm to detect fraudulent transactions.

2 Data Preprocessing

Data preprocessing is crucial for any machine learning task, particularly for K-NN models, which are sensitive to feature scales. For this project, we used Scikit-learn's `RobustScaler`, which scales data using the interquartile range, minimizing the influence of outliers. We split the dataset into training (80%)

and test (20%) sets, fitting the scaler to the training set and transforming both sets accordingly.

```
1 from sklearn.preprocessing import RobustScaler
2 from sklearn.model_selection import train_test_split
3
4 # Splitting the data
5 X_train, X_test, y_train, y_test = train_test_split(X, y, stratify=
    y, test_size=0.2, random_state=42)
6
7 # Scaling the data
8 scaler = RobustScaler()
9 X_train_s = scaler.fit_transform(X_train)
10 X_test_s = scaler.transform(X_test)
```

Listing 1: Data Preprocessing

3 Feature Selection

We excluded features with a correlation close to 0 with the target variable. Specifically, any feature with an absolute correlation value less than 0.1 was removed from the dataset.

```
1 target_correlations = df.corr()['Class'].sort_values(ascending=
    False)
2 selected_features = target_correlations[target_correlations.abs()
    >= 0.1].index.tolist()
3 df_clean = df[selected_features]
```

Listing 2: Feature Selection

4 Model Implementation

We implemented the K-NN algorithm using a custom classifier class that includes methods for fitting the model to training data, predicting labels for test data, and evaluating performance. Euclidean distance was used as the distance metric, and we tested various values of k to identify the optimal number of neighbors. To evaluate performance, we performed 5-fold cross-validation using scikit-learn's `StratifiedKFold`. We averaged the accuracy, precision, and recall for each k value and recorded the results in a dictionary.

```
1 from numba import njit, prange
2
3 @njit(parallel=True, fastmath=True)
4 def euclidean_distance(X_train, X_test):
5     num_train, num_features = X_train.shape
6     num_test = X_test.shape[0]
7     distances = np.empty((num_test, num_train), dtype=np.float64)
8
9     for i in prange(num_test):
10         for j in range(num_train):
11             diff = X_train[j] - X_test[i]
```

```

12         distances[i, j] = np.sqrt(np.sum(diff ** 2))
13
14     return distances
15
16 @njit(parallel=True, fastmath=True)
17 def predict_labels(distances, y_train, k):
18     num_test = distances.shape[0]
19     predictions = np.empty(num_test, dtype=np.int32)
20
21     for i in prange(num_test):
22         neighbors_indices = np.argsort(distances[i])[:k]
23         neighbor_labels = y_train[neighbors_indices]
24         count_1 = np.sum(neighbor_labels == 1)
25         count_0 = np.sum(neighbor_labels == 0)
26         predictions[i] = 1 if count_1 > count_0 else 0
27
28     return predictions
29
30 class KNNClassifier:
31
32     def __init__(self, k=5):
33         self.k = k
34
35     def fit(self, X_train, y_train):
36         self.X_train = X_train
37         self.y_train = y_train
38
39     def predict(self, X_test, batch_size=100):
40         num_samples = X_test.shape[0]
41         predictions = np.empty(num_samples, dtype=np.int32)
42
43         for i in range(0, num_samples, batch_size):
44             end_index = min(i + batch_size, num_samples)
45             batch_X_test = X_test[i:end_index]
46             distances = euclidean_distance(self.X_train,
batch_X_test)
47             batch_predictions = predict_labels(distances, self.
y_train, self.k)
48             predictions[i:end_index] = batch_predictions
49
50         return predictions
51
52     def evaluate(self, X_test, y_test):
53         self._predictions = self.predict(X_test)
54
55         self._accuracy = np.sum(self._predictions == y_test) / len(
y_test)
56
57         # Compute the number of true positives, false positives,
and false negatives
58         true_positives = np.sum((self._predictions == 1) & (y_test
== 1))
59         false_positives = np.sum((self._predictions == 1) & (y_test
== 0))
60         false_negatives = np.sum((self._predictions == 0) & (y_test
== 1))
61

```

```

62         # Compute precision and recall
63         self._precision = true_positives / (true_positives +
false_positives) if (true_positives + false_positives) > 0 else
0
64         self._recall = true_positives / (true_positives +
false_negatives) if (true_positives + false_negatives) > 0 else
0
65
66     def predictions(self):
67         return self._predictions # Return the predictions
68
69     def metrics(self):
70         return np.array([self._accuracy, self._precision, self.
_recall]) # Return the metrics as an array
71
72 def cross_validate_knn(X, y, k_values, n_splits=5):
73     kf = StratifiedKFold(n_splits=n_splits, shuffle=True,
random_state=42)
74     results = []
75
76     for k in k_values:
77         knn = KNNClassifier(k=k)
78         accuracies = np.array([])
79         precisions = np.array([])
80         recalls = np.array([])
81
82         for train_index, val_index in kf.split(X,y):
83             X_train, X_val = X[train_index], X[val_index]
84             y_train, y_val = y[train_index], y[val_index]
85
86             knn.fit(X_train, y_train)
87             knn.evaluate(X_val, y_val)
88             accuracy, precision, recall = knn.metrics()
89
90             accuracies=np.append(accuracies,accuracy)
91             precisions=np.append(precisions,precision)
92             recalls=np.append(recalls,recall)
93
94
95         results=np.append(results,{
96             'k': k,
97             'accuracy': np.mean(accuracies),
98             'precision': np.mean(precisions),
99             'recall': np.mean(recalls)
100         })
101
102
103     print({
104         'k': k,
105         'accuracy': np.mean(accuracies),
106         'precision': np.mean(precisions),
107         'recall': np.mean(recalls)
108     })
109
110     return results

```

Listing 3: Model Implementation

5 Model Evaluation and Results

We tested k values ranging from 5 to 25 in increments of 5. Although each k value performed similarly during hyperparameter testing, $k=5$ yielded the best results. Ultimately, our model with $k=5$ achieved excellent performance on the test set, with 99.95% accuracy, 91.76% precision, and 79.59% recall.

```
1 results = cross_validate_knn(X_train_s, y_train, np.arange(5, 30,  
2     5))  
3 # Example result  
4 {'k': 5, 'accuracy': 0.9995479382913823, 'precision': 0.91764706, '  
    recall': 0.79591837}
```

Listing 4: Model Evaluation

6 Analysis and Discussion

6.1 Feature Selection

The feature selection process of removing features with an absolute correlation of less than 0.1 with the target improved the performance by ignoring features that have a weak relationship with the target variable but have a significant effect on the distance between points. While this primitive technique was useful in the context of our problem, future analysis could benefit from more sophisticated feature selection methods to achieve even better results.

6.2 Robustness and Generalization

We used 5 folds in our cross-validation to ensure that the model's performance could be generalized across different subsets of the data set. Because each subset of the training set performed consistently well, we can be confident the model generalizes to unseen data.

6.3 Model Performance

While the model's high accuracy shows that K-NN can correctly classify transactions most of the time, it certainly does not paint a complete picture. A credit card transaction classification model must have a high precision rate to minimize the number of legitimate transactions flagged as fraudulent, but the most important metric for evaluating the model's performance in this problem is the recall, the proportion of the fraudulent transactions the algorithm classified as fraudulent. In an ideal situation, we like to design a model that classifies every fraudulent transaction as fraudulent, and while a 79.59% recall rate is decent, it shows that this model still has room to improve.

6.4 Limitations and Disadvantages of K-NN

Computational Efficiency: Calculating the distance between points and predicting labels are considerable bottlenecks in evaluating the test set's performance. To help mitigate this, we wrote JIT-compiled functions to increase the efficiency of calculating the distance between points and predicting class labels, but despite significant efforts, K-NN is still computationally expensive for evaluating large data sets. Future analysis could benefit from dividing the training set into regions, so data points in the test set only have to compute the distance of data points that appear in the same region.

Imbalanced Data Set: While we used stratification to ensure similar proportions of fraudulent transactions in different subsets of the data, fraudulent credit card transactions account for less than a fifth of a percent of all points in the data set. Because of the imbalance between classes, the legitimate transactions tend to dominate the voting process for higher values of k leading to lower recall rates.

6.5 Strengths of K-NN

Simplicity: While optimization for computational efficiency is quite difficult, K-NN models are simple to write. There aren't many hyperparameters to test, so hyperparameter tuning doesn't require testing as many combinations.

Interpretability: It's very easy to comprehend how a K-NN model classifies points in the test set.

Incremental Learning: It's easy to add and remove points from the model. In the context of credit card fraud, individuals make transactions every second of the day. When the goal is constantly updating the model with recent transactions, K-NN can easily implement these changes.